1 AGREEMENT AND ATOMIC BROADCAST IN ASYNCHRONOUS
2 NETWORKS

3 TECHNICAL FIELD

4 The present invention relates to a network whose processor nodes exchange information
5 by sending messages in an asynchronous fashion. More particularly, the invention relates
6 to methods and apparatus for achieving agreement among processors, even in the
7 presence of undetected faulty processors, and to a method for reliably broadcasting the
8 messages in an order within an asynchronous point-to-point network. Thus, it is
9 applicable in a wide range of distributed computation systems, reaching from
10 fault-tolerant database systems to intrusion tolerant e-commerce servers.

11 BACKGROUND OF THE INVENTION

12 Distributed systems running in error-prone and adversarial environments have to rely on
13 trusted components. In today's Internet these are typically directory and authorization
14 services, such as the domain name system (DNS), Kerberos, certification authorities, or
15 secure directories. Building such centralized trusted services has turned out to be a
16 valuable design principle for computer security because the trust in them can be leveraged
17 to many, diverse applications that all benefit from centralized management. Often, a
18 trusted service is implemented as the only task of an isolated and physically protected
19 machine.

20 Unfortunately, centralization introduces a single point of failure. Even worse, it is
21 increasingly difficult to protect any single system against the sort of attacks proliferating
22 on the Internet today. One established way for enhancing the fault tolerance of centralized
23 components is to distribute them among a set of servers and to use replication algorithms

CH920000062US1

1   for masking faulty servers or devices. Thus, no single server has to be trusted completely

2   and the overall system derives its integrity from a majority of correct servers.

3   The use of cryptographic methods for maintaining consistent state in a distributed system

4   has a long history and originates with the work of M. Pease, R. Shostak, and L. Lamport,

5   in "Reaching agreement in the presence of faults," Journal of the ACM, vol. 27, pp.

6   228-234, Apr. 1980.

7   The work of M. K. Reiter and K. P. Birman, "How to securely replicate services," ACM

8   Transactions on Programming Languages and Systems, vol. 16, pp. 986-1009, May 1994

9   introduces secure state machine replication in a Byzantine environment and a broadcast

10   protocol based on threshold cryptography that maintains causality among the requests.

11   Since no robust threshold-cryptographic schemes and secure atomic broadcast protocols

12   were not known at that time, no fully robust systems for an asynchronous environment

13   with malicious faults could be designed.

14   Subsequent work by Reiter in "Distributing trust with the Rampart toolkit,"

15   Communications of the ACM, vol. 39, pp. 71-74, Apr. 1996 assumes a model which

16   implements atomic broadcast on top of a group membership protocol that dynamically

17   removes apparently faulty servers from the set.

18   M. Castro and B. Liskov in "Practical Byzantine fault tolerance," in Proc. Third Symp.

19   Operating Systems Design and Implementation, 1999 present a practical algorithm for

20   distributed service replication that is fast if no failures occur. It requires no explicit

21   time-out values, but assumes that message transmission delays do not grow faster than

22   some predetermined function for an indefinite duration. Since this protocol is

23   deterministic, it can be blocked by a Byzantine adversary (i.e., violating liveness). In

24   contrast, an approach based on a probabilistic agreement protocol satisfying both

25   conditions would be a better approach.

1 The *Total* family of algorithms for total ordering by L. E. Moser and P. M. Melliar-Smith,
2 "Byzantine-resistant total ordering algorithms," Information and Computation, vol. 150,
3 pp. 75-111, 1999 implements atomic broadcast in a Byzantine environment, but only
4 assuming a benign network scheduler with some specific probabilistic fairness
5 guarantees. Although this may be realistic in highly connected environments with
6 separate physical connections between all machines, it seems not appropriate for arbitrary
7 Internet settings.

8 K. P. Kihlstrom, L. E. Moser, and P. M. Melliar-Smith, "The SecureRing protocols for
9 securing group communication," in Proc. 31st Hawaii International Conference on
10 System Sciences, pp. 317-326, IEEE, Jan. 1998 and the work of A. Doudou, B.
11 Garbinato, and R. Guerraoui, "Abstractions for devising Byzantine-resilient state machine
12 replication," in Proc. 19th Symposium on Reliable Distributed Systems (SRDS 2000), pp.
13 144-152, 2000 are two examples of atomic broadcast protocols that rely on failure
14 detectors in the Byzantine model. They encapsulate all time-dependent aspects and
15 obvious misbehavior of a party in the abstract notion of a failure detector and permit
16 clean, deterministic protocols. But failure detectors are not well understood in Byzantine
17 environments.

18 US Patent 4,644,542 describes a method for reliably broadcasting information in a
19 point-to-point network of processors in the presence of component faults provided that
20 the network remains connected using only an exchange of messages. The method
21 possesses the properties that every message broadcast by a fault-free processor is accepted
22 exactly once by all fault-free processors within a bounded time, that every message
23 broadcast is either accepted by all fault-free processors or none of them, and that all
24 messages accepted by fault-free processors are accepted in the same order by all those
25 processors. The method is based on a diffusion technique for broadcasting information
26 and on special message validity tests for tolerating any number of component failures up
27 to network partitioning or successful forgery.

CH920000062US1

1    US Patent 5,598,529 discloses a computer system resilient to a wide class of failures

2    within a synchronized network. It includes a consensus protocol, a broadcast protocol and

3    a fault tolerant computer system created by using the two protocols together in

4    combination. The protocols are subject to certain validity conditions. The system in the

5    state of consensus is guaranteed to have all non-faulty processors in agreement as to what

6    action the system should take. The system and protocols can tolerate up to $t$ processor

7    failures out of $3t+1$ or more processors, but requires as well as the before mentioned

8    method timing guarantees and is therefore not suitable for asynchronous networks.

9    Fault-tolerant systems use computer programs called protocols to ensure that the systems

10    will operate properly even if there are individual processor failures.

11    A fault-tolerant consensus protocol enables each processor or party to propose an action

12    (via a signal) that is required to be coordinated with all other processors in the system. A

13    fault-tolerant consensus protocol has as its purpose the reaching of a "consensus" on a

14    common action (e.g., turning a switch off or on) to be taken by all non-faulty processors

15    and ultimately the system. Consensus protocols are necessary because processors may

16    send signals to only a single other processor at a time and a processor failure can cause

17    two processors to disagree on the signal sent by a third failed processor. In spite of these

18    difficulties, a fault-tolerant consensus protocol ensures that all non-faulty processors

19    agree on a common action.

20    To reach consensus, consensus protocols first enable each processor or participating

21    network device to propose an action (via a signal) that is later to be coordinated by all the

22    processors or participating network devices in the system. The system then goes through

23    the steps of the consensus protocol. After completing the consensus protocol steps, the

24    common action of the consensus is determined.

25    A safe architecture for distributing trusted services among a set of servers is desired that

26    guarantees availability and integrity of the services despite some servers being under

1   control of an attacker or failing in arbitrary malicious ways. The architecture should be

2   characterized by a static set of servers and completely asynchronous point-to-point

3   communication. Trusted applications can only be achieved by an efficient and provably

4   secure agreement and broadcast protocol.

5   SUMMARY OF THE INVENTION

6   One aspect of the invention provides a machine-implementable method for achieving

7   multi-valued Byzantine agreement among $n$ participating network devices, also called

8   processors or parties. The method achieves agreement and consensus among the $n$

9   participating network devices in an asynchronous network for deciding on a common

10  value $v$. The common value $v$ is validated by a justification $p$, whereby the common value

11  $v$ and the justification $p$ together satisfy a predetermined predicate $Q$. The agreement

12  arises out of a series of messages being sent and received by each participating network

13  device with up to a number $t$ of faulty devices. The method turns out to be practical and

14  also theoretically nearly optimal in the sense that it withstands the maximum number of

15  corrupted devices $t < n/3$ and uses a nearly optimal number of messages whereby the total

16  bit-length of these messages is also nearly optimal.

17  Another aspect of the invention provides a method for reliably broadcasting messages in

18  an order within the asynchronous network comprising $n$ participating network devices.

19  This ordered broadcast protocol, also referred to as atomic broadcast protocol, tolerates a

20  number $t$ of less than $n/3$ faulty participating network devices. Client requests to a service

21  are delivered by the atomic broadcast protocol that imposes a total order on all requests

22  and guarantees that the participating network devices perform the same sequence of

23  operations. The atomic broadcast can be realized by use of a randomized protocol to solve

24  Byzantine agreement, such as the multi-valued Byzantine agreement as mentioned above.

1    In another aspect of the invention a secure architecture for distributing trusted services

2    among a set of participating network device is provided by use of the multi-valued

3    Byzantine agreement and the atomic broadcast protocol that guarantees availability and

4    integrity (or equivalently, liveness and safety) of the services despite some participating

5    network devices being under control of an attacker or failing in arbitrary malicious ways.

6    The architecture is characterized by a static set of participating network device,

7    completely asynchronous point-to-point communication, and the use of modern

8    cryptographic techniques. Trusted applications can be realized by deterministic state

9    machines replicated on all servers or participating network devices and initialized to the

10    same state.

11    In the asynchronous model there are no timing assumptions in the design of any protocol.

12    Asynchronous protocols are attractive because in a synchronous system, the designer has

13    to specify time-out values, which is very difficult when protecting against arbitrary

14    failures that may be caused by a malicious attacker. It is usually much easier for an

15    intruder to block communication with a server than to subvert it. Prudent security

16    engineering also gives the adversary full access to all specifications, including time-outs,

17    and excludes only cryptographic keys. Such an adversary may simply delay the

18    communication to a server for a little longer than the time-out and the server appears

19    faulty to the remaining system. Therefore the presented asynchronous protocols are

20    designed to resist all timing attacks.

21    The basic idea of the multi-valued Byzantine agreement is that every participating

22    network device proposes its value as a candidate value for the final result. One

23    participating network device whose proposal satisfies the validation predicate $Q$ is then

24    selected in a sequence of binary Byzantine agreement protocols and this value becomes

25    the final decision value.

1     In general, the method for achieving multi-valued Byzantine agreement among $n$

2     participating network devices comprises: i) echoing a proposal, ii) an agreement loop, and

3     iii) delivering the chosen proposal. More precisely, the multi-valued Byzantine agreement

4     comprises the detailed steps: i) broadcasting to the participating network devices an echo

5     message comprising a proposed value $w$ and a proposed justification $\pi$ by using verifiable

6     and consistent broadcast; ii) receiving $n-t$ echo messages comprising candidate values $w_i$

7     and candidate justifications $\pi_i$ satisfying the predicate $Q$, and repeating the following

8     steps 1) to 3) for each participating network device as a candidate device represented by a

9     candidate device identifier a according to an order: 1) broadcasting to all participating

10.   network devices a vote message comprising the candidate device identifier a, and either a

11    first agree-value Y together with the candidate value $w_a$ and the candidate justification $\pi_a$,

12    or a second agree-value N, 2) receiving vote messages and counting up to $n-t$ vote

13    messages including the second agree-value N or the first agree-value Y, the candidate

14    value $w_a$, and the candidate justifications $\pi_a$ satisfying the predicate $Q$, and 3) performing

15    a Byzantine agreement to determine whether the candidate device has sent the candidate

16    value $w_a$ and the candidate justification $\pi_a$ satisfying the predicate $Q$, iii) in response to

17    the result of the Byzantine agreement, deciding the common value $v$ proposed as the

18    candidate value $w$ and the justification $p$ proposed as the candidate justification $\pi$ of an

19    agreed candidate device.

20    The multi-valued Byzantine agreement protocol may further comprise in step ii)

21    broadcasting a commit message comprising sender identities of received echo messages,

22    receiving commit messages, and selecting randomly the order in which the candidate

23    device is considered by opening at least one cryptographic common coin. By doing so, the

24    counted vote messages in step 2) can be restricted to those that are consistent with the

25    received commit messages. This has the advantage that the protocol runs in a constant

26    expected number of rounds and therefore terminates faster. In particular, the agreement

1     loop is executed a number of times which is independent of the number of participating

2     network devices.

3     The step of opening at least one cryptographic common coin may comprise using a

4     distributed coin-tossing protocol, which has the advantage that a distributed protocol can

5     be used and no centralized authority is required.

6     The verifiable and consistent broadcast in step i) may comprise a certified delivery of the

7     broadcasted echo message within a sent message and whereby the agreement is reached

8     on the content of the broadcasted echo message. This shows the advantage that certified

9     delivery occurs within a single message.

10     Furthermore, the verifiable and consistent broadcast in step i) may comprise an exchange

11     of signed messages between the participating network devices.

12     The verifiable and consistent broadcast in step i) may comprise the use of threshold

13     signatures. By using them, the size of the used messages can be reduced. A suitable

14     threshold signature scheme has been provided by V. Shoup and published in the article

15     "Practical threshold signatures", in Advances in Cryptology: EUROCRYPT 2000 (B.

16     Preneel, ed.), vol. 1087 of Lecture Notes in Computer Science, pp. 207-220, Springer,

17     2000. This article is incorporated herein by means of reference.

18     The multi-valued Byzantine agreement protocol may further comprise in steps 1) and 2)

19     in step ii) broadcasting the candidate value $w_a$ and the candidate justification $\pi_a$ only upon

20     a request. This proofs advantageously, because it reduces the number and the size of the

21     messages.

22     The multi-valued Byzantine agreement protocol may further comprise in step ii) that the

23     participating network devices are voting on several candidate devices simultaneously.

24     This shows the advantage that if the number $n$ of participating network devices is small

25     then the number of rounds can be kept small as well.

CH920000062US1

1 Within the atomic broadcast protocol each participating network device stores a queue $q$

2 and a log file, hereafter short log $d$. The atomic broadcast protocol operates in rounds,

3 whereby each round comprises the following steps: i) responsive to a message broadcast

4 request comprising a message value $m$ performing the step of: appending the message

5 value $m$ to the queue $q$ unless the log $d$ or the queue $q$ comprises the message value $m$, ii)

6 deriving a signature $\sigma$ on the queue $q$, iii) broadcasting to all participating network

7 devices a queue message comprising the queue $q$ and the signature $\sigma$, iv) receiving a

8 number $c$ of at least $t+1$ queue messages comprising $c$ proposed queues $q_i$ and proposed

9 signatures $\sigma_i$, v) storing the proposed queues $q_i$ in a queue vector QV and the proposed

10 signatures $\sigma_i$ in a signature vector SV, vi) proposing the queue vector QV for Byzantine

11 agreement validated by the signature vector SV and performing a method for achieving

12 agreement on a common value being validated by a justification $p$ together satisfying a

13 predetermined predicate $Q$ by validating the queue vector QV and the signature vector SV

14 through a determined predicate Q asserting that the signature vector SV comprises $c$ valid

15 signature entries of distinct participating network devices on entries of the queue vector

16 QV, vii) preparing in response to the result of the Byzantine agreement an ordered list L

17 of unique message values out of the entries of the decided queue vector DQV, viii)

18 accepting the unique message values in the ordered list L in the sequence of the ordered

19 list L, ix) appending the accepted unique message values to the log $d$.

20 The atomic broadcast protocol may further comprises in step iv) appending an unknown

21 message value found in a received queue message to the queue $q$ unless the log $d$ or the

22 queue $q$ comprises the unknown message value. This has the advantage that the unknown

23 message is delivered faster by the system.

24 **Glossary**

25 The following are helpful definitions to aid in the understanding of the description.

CH920000062US1

1     **Validated Byzantine Agreement**: A protocol solves validated Byzantine agreement with

2     predicate $Q$ if it satisfies the following conditions except with negligible probability:

3         **External Validity**: Any honest participating network device that terminates decides

4         on a common value $v$ validated by a justification $\pi$ such that $Q(v, \pi)$ holds.

5         **Agreement**: If some honest participating network device decides on the common

6         value $v$, then any honest participating network device that terminates decides on the

7         common value $v$.

8         **Termination**: Each honest participating network device eventually terminates the

9         validated Byzantine agreement protocol.

10     Thus, honest participating network devices may propose all different values and the

11     decision value may have been proposed by a corrupted participating network device, as

12     long as honest participating network devices obtain the corresponding validation during

13     the protocol. The agreement and termination are the same as in the definition of ordinary,

14     binary Byzantine agreement. Another variation of the validity condition is that an

15     application may prefer one decision value over others. Such an agreement protocol may

16     be biased and always output the preferred value in cases where other values would have

17     been valid as well. For binary validated agreement, the method of the invention will apply

18     a protocol that is biased towards 1 or Yes. Its purpose is to detect whether there is a

19     validation for 1, so it suffices to guarantee termination with output 1 if $t+1$ honest

20     participating network devices know the corresponding information at the outset. A binary

21     validated Byzantine agreement protocol biased towards 1 is a protocol for validated

22     Byzantine agreement on values in $\{0, 1\}$ such that the following condition holds:

1     **Biased Validity**: If at least $t+1$ honest participating network devices propose 1 validated

2     by $\pi$ such that $Q(1, \pi)$ holds, then any honest participating network device that terminates

3     decides 1.

4     **Reliable Broadcast**: In a reliable broadcast each message is delivered with an indication

5     of its sender. A reliable broadcast protocol satisfies the following properties:

6        **Validity**: If an honest participating network device broadcasts a message $m$, then all

7        honest participating network devices eventually deliver $m$.

8        **Agreement**: If some honest participating network device has delivered $m$, then

9        eventually all honest participating network devices deliver $m$.

10        **Integrity**: Each honest participating network device delivers $m$ at most once, and if

11        the indicated sending participating network device is honest, then $m$ was previously

12        sent by the indicated participating network device.

13     **Verifiable Broadcast**: A reliable broadcast protocol is called verifiable if the following

14     holds, except with negligible probability: When an honest participating network device

15     has delivered a message $m$, then it can produce a single protocol message $M$ that it may

16     send to other participating network devices such that any other honest participating

17     network device will deliver $m$ upon receiving $M$ (provided the other participating network

18     device has not already delivered $m$).

19     $M$ is the message that completes the verifiable broadcast. This notion implies that there is

20     a predicate that the receiving participating network device can apply to an arbitrary bit

21     string for checking if it constitutes a message that completes a verifiable broadcast.

22     **Consistent Broadcast**: A protocol for consistent broadcast satisfies all properties of a

23     reliable broadcast, except that the agreement property is replaced by the following:

**Consistency**: If some honest participating network device delivers $m$ and another honest participating network device delivers $m'$, then $m = m'$.

Thus, consistent broadcast makes no provisions that two participating network devices deliver the message, but maintains agreement among the actually delivered messages with the same senders and sequence numbers.

**Atomic Broadcast**: Atomic broadcast is an ordered reliable broadcast. A protocol for atomic broadcast is a protocol for reliable broadcast that also satisfies the following property:

**Total Order**: If a first honest participating network device delivers $m$ and $m'$, and a second honest participating network device delivers $m$ and $m'$, then the first honest participating network device delivers $m$ before $m'$ if and only if the second honest participating network device delivers $m$ before $m'$.

It is known that protocols for atomic broadcast are considerably more expensive than those for reliable broadcast. This is because even with crash faults, atomic broadcast is equivalent to consensus and cannot be solved by deterministic protocols in an asynchronous network.

**Digital Signature**: The invention uses unforgeable, transferable digital signatures to allow parties or participating network devices to justify their messages, which limits the possibilities of an adversary.

**Threshold cryptography**: Threshold cryptographic schemes are non-trivial extensions of the classical concept of secret sharing in cryptography. Secret sharing allows a group of $n$ participating network devices or parties to share a secret such that $t$ or fewer of them have no information about it, but $t+1$ or more can uniquely reconstruct it. However, one cannot simply share the secret key of a cryptosystem and reconstruct it for decrypting a message

CH920000062US1

1 because as soon as a single corrupted party knows the key, the cryptosystem becomes

2 completely insecure and unusable.

3 **Threshold Signature**: In a threshold signature scheme, each participating network device

4 holds a share of the secret signing key and may generate shares of signatures on

5 individual messages upon request. The validity of a signature share can be verified for

6 each participating network device. From $t+1$ valid signature shares, one can generate a

7 digital signature on the message that can later be verified using the single, publicly known

8 signature verification key. In a secure threshold signature scheme, it is infeasible for a

9 computationally bounded adversary to produce $t+1$ valid signature shares that cannot be

10 combined to a valid signature (robustness), and to output a valid signature on a message

11 for which no honest participating network device generated a signature share (no forgery).

12

13 **Threshold Coin-Tossing**: A threshold coin-tossing protocol provides arbitrarily many

14 unpredictable random bits. An efficient implementation of a threshold coin-tossing

15 scheme is part of the randomized Byzantine agreement protocol of C. Cachin, K.

16 Kursawe, and V. Shoup, "Random oracles in Constantinople: Practical asynchronous

17 Byzantine agreement using cryptography," in Proc. 19th ACM Symposium on Principles

18 of Distributed Computing (PODC), pp. 123-132, 2000. It can be used to ensure

19 termination of the agreement protocol within an expected constant number of rounds.

20 **Hybrid Failures**

21 The method for achieving agreement among $n$ participating network devices in an

22 asynchronous network for deciding on a common value $v$ and the method for reliably

23 broadcasting messages in an order can distinguish between several different ways in

24 which a network device can fail. This could for example be

1      **Byzantine Failures BF:** If a byzantine failure BF occurs, the adversary has taken full

2      control over the corresponding machine. All secrets this machine has are handed over

3      to the adversary, who now controls its entire behavior.

4      **Crash Failures CF:** A crash failure CF simply means that the corresponding

5      machine stops working. This could happen anytime, i.e., even in the middle of a

6      broadcast or while sending a message. It is assumed that there is no mechanism other

7      parties can reliably detect such a crash.

8      **Link Failures LF:** A link failure LF occurs when not a party, but an interconnecting

9      link becomes faulty. As the link has no access to authentication keys, it is easy to

10     prevent it from modifying or inserting messages. A faulty link could however delete

11     messages, and it might completely disconnect two parties.

12     **Adversary structure**

13     An adversary structure $T$ is a set of sets (coalitions) of parties whose corruption the

14     system should be tolerated. Let $M$ be the set of all participating network devices. An

15     adversary structure is called

16     $Q^2$, if no two coalitions $N_1, N_2 \in T$ satisfy $N_1 \cup N_2 = M$.

17     $Q^3$, if no three coalitions $N_1, N_2, N_3 \in T$ satisfy $N_1 \cup N_2 \cup N_3 = M$.

18     $Q^{2+3}$ with respect to CF and BF, if for all $c_1, c_2 \in$ CF and all $b_1, b_2, b_3 \in$ BF,

19     $M \setminus \{b_1 \cup b_2 \cup b_3 \cup c_1 \cup c_2\} \supsetneq \varnothing$;

20     A $Q^2$ adversary structure is sufficient to solve byzantine agreement if only crash failures

21     CF occur. $Q^3$ is applied in the byzantine case, where only byzantine failures BF occur,

22     while $Q^{2+3}$ is the generalization for the hybrid crash-byzantine failure case.

23

1 DESCRIPTION OF THE DRAWINGS

2 The invention is described in detail below with reference to the following schematic

3 drawings.

4 **FIG. 1** shows a typical asynchronous network with multiple participating network

5 devices.

6 **FIG. 2** shows an example of a schematic diagram of an agreement protocol

7 according to the present invention.

8 **FIG. 3a** shows an extension of the embodiment of the agreement protocol of FIG.

9 2.

10 **FIG. 3b** shows a further extension of the embodiment of the agreement protocol of

11 FIG. 2.

12 **FIG. 4** shows a an example of a schematic illustration of an atomic Broadcast

13 protocol according to the present invention.

14 DETAILED DESCRIPTION OF THE INVENTION

15 The present invention relates to a network whose processor nodes exchange information

16 by sending messages in an asynchronous fashion. More particularly, the invention relates

17 to a methods and apparatus for achieving agreement among the processors, even in the

18 presence of undetected faulty processors, and to a method for reliably broadcasting the

19 messages in an order within an asynchronous point-to-point network.

1 With general reference to the figures, the essential features of a method for achieving

2 agreement among $n$ participating network devices in an asynchronous network for

3 deciding on a common value $v$ and a method for reliably broadcasting messages in an

4 order are described in more detail below.

5 At first, some basics, in accordance with the present invention, are addressed.

6 *Validated Byzantine Agreement*

7 The standard notion of Byzantine agreement implements a binary decision and can

8 guarantee a particular outcome only if all honest participating network devices propose

9 the same value. A weaker validity condition is herewith introduced, called *external*

10 *validity*, which relaxes the validity condition and generalizes to decisions on a common

11 value $v$ from an arbitrarily large set. It requires that the decided value $v$ satisfies a global

12 predicate $Q$ that is determined by the particular application and known to all participating

13 network devices. Each participating network device adds some validation data to the

14 proposed value, which serves as the proof for its validity. Typically, this comprises a

15 digital signature that can be verified by all participating network devices. The agreement

16 protocol then returns to a caller not only the decision value, but also the corresponding

17 validation data. The caller might need this information if it did not know it before.

18 The validated Byzantine agreement generalizes the primitive of agreement on a core set

19 and the notion of interactive consistency, as described by M. J. Fischer, "The consensus

20 problem in unreliable distributed systems (a brief survey)," in Foundations of

21 Computation Theory (M. Karpinsky, ed.), vol. 158 of Lecture Notes in Computer

22 Science, Springer, 1983, to the Byzantine model, which uses agreement on a vector of $n$

23 values, one from each participating network device.

24 *Verifiable Broadcast*

CH920000062US1

1    One participating network device that has delivered a payload message using reliable

2    broadcast may want to inform another respective participating network device about this.

3    Such information might be useful to the respective participating network device if it has

4    not yet delivered the message, but can exploit this knowledge somehow, in particular

5    since the respective participating network device is guaranteed to deliver the same

6    message by the agreement property in reliable broadcast. In a standard reliable broadcast,

7    however, this knowledge cannot be transferred in a verifiable way. Therefore, this

8    property of a broadcast protocol is identified here because it is useful in the application,

9    and called *verifiability*. Informally, it means this: when the respective participating

10    network device claims that it is not yet in a state to deliver a particular payload message

11    $m$, then the one participating network device can reply with a single protocol message and

12    when the respective participating network device processes this, it will deliver $m$

13    immediately and terminate the corresponding broadcast.

14    *Reliable Broadcast*

15    Protocols for reliable broadcast in an asynchronous network with Byzantine faults are

16    well known in the art, where the reliable broadcast is also known as the problem of the

17    *Byzantine generals*, as described in the article by L. Lamport, R. Shostak, and M. Pease,

18    "The Byzantine generals problem," ACM Transactions on Programming Languages and

19    Systems, vol. 4, pp. 382-401, July 1982. An applicable protocol has been described by G.

20    Bracha in the article "An asynchronous $\lfloor \frac{n-1}{3} \rfloor$-resilient consensus protocol," in Proc. 3rd

21    ACM Symposium on Principles of Distributed Computing (PODC), pp. 154-162, 1984.

22    *Consistent Broadcast*

23    Several protocols for consistent broadcast have been proposed by M. Reiter, "Secure

24    agreement protocols: Reliable and atomic group multicast in Rampart," in Proc. 2nd

25    ACM Conference on Computer and Communications Security, 1994. The agreement

1    property of reliable broadcast is rather expensive to satisfy; it is the main reason why

2    most protocols for reliable broadcast need on the order of $n^2$ messages. For some

3    applications, however, agreement is not necessary and can be ensured by other means, as

4    long as integrity is satisfied. The resulting notion is called consistent broadcast.

5    *Cryptography*

6    Cryptographic techniques such as public-key encryption schemes and digital signatures

7    are useful for many existing secure services. For distributing trusted services, distributed

8    variants of them from threshold cryptography can be applied.

9    Threshold-cryptographic protocols have been used for secure service replication before,

10    e.g., as described by M. K. Reiter and K. P. Birman, "How to securely replicate services,"

11    ACM Transactions on Programming Languages and Systems, vol. 16, pp. 986-1009, May

12    1994. However, a major complication for adopting threshold cryptography to an

13    asynchronous distributed system is that many early protocols are not robust and that most

14    protocols rely heavily on synchronous broadcast channels. Only very recently,

15    non-interactive schemes have been developed that satisfy the appropriate notions of

16    security, such as the threshold cryptosystem of V. Shoup and R. Gennaro, "Securing

17    threshold cryptosystems against chosen ciphertext attack," in Advances in Cryptology:

18    EUROCRYPT '98 (K. Nyberg, ed.), vol. 1403 of Lecture Notes in Computer Science,

19    Springer, 1998 and the threshold signature scheme of V. Shoup, "Practical threshold

20    signatures," in Advances in Cryptology: EUROCRYPT 2000 (B. Preneel, ed.), vol. 1087

21    of Lecture Notes in Computer Science, pp. 207-220, Springer, 2000.

22    *Non-Interactive Threshold Signatures*

23    A useful tool for the atomic broadcast protocol are non-interactive threshold signatures.

24    More precisely, a dual-threshold variations can be applied as introduced by C. Cachin, K.

25    Kursawe, and V. Shoup, in "Random oracles in Constantinople: Practical asynchronous

1    Byzantine agreement using cryptography," in Proc. 19th ACM Symposium on Principles

2    of Distributed Computing (PODC), pp. 123-132, 2000. The basic idea of a dual-threshold

3    signature scheme is that there are $n$ participating network devices, $t$ of which may be

4    corrupted. The participating network devices hold shares of the secret key of a signature

5    scheme and may generate shares of signatures on individual messages. The condition is

6    that $\kappa$ signature shares are necessary and sufficient to construct a signature, where $t < \kappa \le$

7    $n - t$.

8    More precisely, a non-interactive $(n, \kappa, t)$-dual-threshold signature scheme comprises the

9    following parts:

10    - A key generation algorithm with input parameters $k$, $n$, $\kappa$ and $t$. It outputs the public key

11    of the scheme, a private key share for each participating network device, and a local veri-

12    fication key for each participating network device.

13    - A signing algorithm with inputs a message, the public key and a private key share. It

14    outputs a signature share on the submitted message.

15    - A share verification algorithm with inputs a message, a signature share on that message

16    from a participating network device, along with the global public key and the local veri-

17    fication key of the respective participating network device. It determines if the signature

18    share is valid.

19    - A share combining algorithm that takes as input a message and $\kappa$ valid signature shares

20    on the message, along with the public key and the verification keys, and outputs a valid

21    signature on the message.

22    - A signature verification algorithm that takes as input a message and a signature

23    (generated by the share-combining algorithm), along with the public key, and determines

24    if the signature is valid.

1  During initialization, a dealer runs the key generation algorithm and gives each

2  participating network device the public key, all local verification keys, and its private key

3  share. The adversary may submit signing requests to the honest participating network

4  devices for messages of its choice. Upon receiving such a request, a participating network

5  device computes a signature share for the given message using its private key share.

6  Given $\kappa$ valid signature shares from distinct participating network devices on the same

7  message, they may be combined into a signature on the message.

8  The two basic security requirements are robustness and non-forgeability. Robustness

9  means that it is computationally infeasible for an adversary to produce $\kappa$ valid signature

10  shares such that the output of the share combining algorithm is not a valid signature.

11  Non-forgeability means that it is computationally infeasible for the adversary to output a

12  valid signature on a message that was submitted as a signing request to less than $\kappa - t$

13  honest participating network devices.

14  A practical scheme that satisfies these definitions in the random oracle model was

15  proposed by V. Shoup, in "Practical threshold signatures," in Advances in Cryptology:

16  EUROCRYPT 2000 (B. Preneel, ed.), vol. 1087 of Lecture Notes in Computer Science,

17  pp. 207-220, Springer, 2000.

18  *Threshold Coin-Tossing*

19  A distributed $(n, t + 1)$-threshold coin-tossing scheme can be applied. The basic idea is

20  the same as for the other threshold primitives, but here the participating network devices

21  hold shares of a pseudorandom function $F$. It maps a bit string $N$, the name of a coin, to

22  its value $F(N) \in \{0, 1\}^{k''}$, whereby a generalized coin is used that produces $k''$ random bits

23  simultaneously. The participating network devices may generate shares of a coin and $t+1$

24  shares of the same coin are both necessary and sufficient to construct the value of that

25  coin. The generation and verification of coin shares are non-interactive.

1    During initialization the dealer generates a global verification key, a local verification key

2    for each participating network device, and a secret key share for each participating

3    network device. The initial state information for each participating network device

4    comprises its secret key share and all verification keys. The secret keys implicitly define

5    the function $F$ mapping names to $k''$-bit strings.

6    After the initialization phase, the adversary submits reveal requests to the honest parties

7    for coins of his choice. Upon receiving such a request, a participating network device

8    outputs a coin share for the given coin computed from its secret key.

9    The coin-tossing scheme also specifies two algorithms:

10    - A share verification algorithm takes as input the name of a coin, a share of this coin

11    from a participating network device, along with the global verification key and the veri-

12    fication key of the respective participating network device, and determines if the coin

13    share is valid.

14    - A share combining algorithm takes as input a name $N$ of a coin and $t+1$ valid shares of

15    $N$, along with the verification keys, and outputs $F(N)$.

16    *Verifiable and Consistent Broadcast*

17    A protocol that implements verifiable and consistent broadcast uses a non-interactive $(n,$

18    $\lceil \frac{n+t+1}{2} \rceil$, $t)$-dual-threshold signature scheme with verifiable shares according to the section

19    above: Non-Interactive Threshold Signatures. All messages are authenticated.

20    The protocol is based on the "echo broadcast" of M. Reiter, "Secure agreement protocols:

21    Reliable and atomic group multicast in Rampart," in Proc. 2nd ACM Conference on

22    Computer and Communications Security, 1994, but uses a threshold signature to decrease

23    the bit complexity. The idea behind it is that the sender broadcasts the message to all

1     participating network devices and hopes for $\lceil \frac{n+t+1}{2} \rceil$ participating network devices to sign

2     it as "witnesses" to guarantee integrity. The signature shares are then collected by the

3     sender, combined to a signature on the message, and relayed to all participating network

4     devices. After receiving the message together with a valid signature, one participating

5     network device delivers it immediately. Because one participating network device may

6     forward the message and the signature to other participating network devices, the protocol

7     is verifiable.


8     *Protocol for Binary Byzantine Agreement*


9     Binary asynchronous Byzantine agreement protocols can be adapted to external validity.

10     For example, in the protocol of C. Cachin, K. Kursawe, and V. Shoup, "Random oracles

11     in Constantinople: Practical asynchronous Byzantine agreement using cryptography," in

12     Proc. 19th ACM Symposium on Principles of Distributed Computing (PODC), pp.

13     123-132, 2000 one has to "justify" the pre-votes of round 1 with a valid $\pi$. The logic of

14     the protocol guarantees that either a decision is reached immediately or the validations for

15     0 and for 1 are seen by all participating network devices in the first two rounds.

16     Furthermore, the protocol can be biased towards 1 by modifying the coin such that it

17     always outputs 1 in the first round.


18     Turning now to Fig. 1 which shows an example of a common computer system 8. It

19     consists of four participating network devices A, B, C, D, which are connected via

20     communication lines (1 through 4 and 5) to a network. The system, where the

21     multi-valued Byzantine Agreement can be achieved, has one faulty device, which is

22     designated by the "X" in participating network device D. Each participating network

23     device A, B, C, D may be any type of computer device known in the art from a computer

24     on a chip or a wearable computer to a large computer system. The communication lines

25     can be any communication means commonly known to transmit data or messages from

26     one participating network device A, B, C, D to another. For instance, the communication

27     lines may be either single, bi-directional communication lines 5 between each pair of

CH920000062US1

1  participating network devices A, B, C, D or one unidirectional line in each direction

2  between each pair of participating network devices A, B, C, D. Such a computer system 8

3  and communication lines 5 are well known in the art. In the case where a participating

4  network device A, B, C, D sends information to itself, an equivalent result could be

5  achieved by merely moving data within the participating network device and not sending

6  it over a communication line to itself. The common computer system 8 is shown to

7  facilitate the description of the following multi-valued Byzantine Agreement protocol.

8  The following steps indicate a general method for achieving multi-valued Byzantine

9  agreement in the validated form, whereby a series of messages being sent and received by

10  each participating network device A, B, C, D. Every participating network device

11  proposes its value as a candidate value $w$ for the final result. One participating network

12  device whose proposal satisfies the validation predicate $Q$ is then selected in a sequence

13  of binary Byzantine agreement protocols and this value becomes the final decision value

14  $v$.

15  In general, the method for achieving multi-valued Byzantine Agreement among $n$

16  participating network devices A, B, C, D comprises the following general steps:

17  i) Echoing a proposal: Each participating network device A, B, C, D sends a proposed

18  value $w$ that it proposes to all other participating network device A, B, C, D using veri-

19  fiable and consistent broadcast. This ensures that all honest participating network devices

20  A, B, C obtain the same proposal value $w$ for any particular participating network device,

21  even if the sender is corrupted. Then, each participating network device A, B, C, D waits

22  until it has received $n - t$ proposals satisfying the predicate $Q$ before entering the

23  following agreement loop.

24  ii) Agreement loop: One participating network device A, B, C, D is chosen after another,

25  according to a fixed permutation $\Pi$ of $\{1, ..., n\}$. Let "a" denote the index of the

26  participating network device A, B, C, D selected in the current round, also referred to as

1     candidate device identifier. For each participating network device A, B, C, D the

2     following steps are executed:

3     1) broadcasting to all participating network devices A, B, C, D a vote message

4     comprising the candidate device identifier a, and either a first agree-value Y together with

5     the candidate value $w_a$ and the candidate justification $\pi_a$, or a second agree-value N,

6     2) receiving vote messages and counting up to $n\text{-}t$ vote messages including the second

7     agree-value N or the first agree-value Y, the candidate value $w_a$, and the candidate

8     justifications $\pi_a$ satisfying the predicate $Q$, and

9     3) performing a validated binary Byzantine agreement biased to 1 or Yes to determine

10    whether the candidate device has sent the candidate value $w_a$ and the candidate

11    justification $\pi_a$ satisfying the predicate $Q$, i.e. voting 1 or Yes if the participating network

12    device A, B, C, D has received a valid proposal from the candidate device and validating

13    this by the protocol message that completes the verifiable broadcast of the candidate

14    device's proposal, otherwise voting 0 or No if the candidate device has received $n\text{-}t$ vote

15    messages comprising 0 or No, and if the validated binary Byzantine agreement decides 1

16    or Yes, exit from the Agreement loop and define the current candidate device as an

17    agreed candidate device.

18    iii) Delivering the chosen proposal: In response to the result of the Byzantine agreement,

19    the common value $v$ is decided which was proposed as the candidate value $w$ and the

20    justification $p$ that was proposed as the candidate justification $\pi$ of said agreed candidate

21    device.

22    In general, a message from one participating network device, e.g., participating network

23    device A to another participating network device, e.g., participating network device B,

24    has the form

1    *ID*, A, B, payload

2    whereby *ID* indicates a unique transaction-identifier. Every instance of a protocol is

3    associated with such a transaction-identifier. A indicates the sender and B the receiver of

4    the message sent in the network, as depicted in Fig. 1.

5    FIG. 2 shows the flow of a multi-valued Byzantine Agreement protocol 100 for an

6    implementation. Each participating network device A, B, C, D performs the steps as

7    indicated by the boxes 10 to 90.

8    At first, as indicated with box 10, the participating network devices A, B, C, D broadcasts

9    an echo message comprising a proposed value $w$ and a proposed justification $\pi$ by using

10   verifiable and consistent broadcast. Then, each participating network device A, B, C, D

11   receives $n\text{-}t$ echo messages comprising candidate values $w_1$, $w_2$, $w_3$ and candidate

12   justifications $\pi_1$, $\pi_2$, $\pi_3$ satisfying the predicate $Q$ as indicated with box 20. The steps in

13   box 30 are an extension of the current agreement and are described with reference to Fig.

14   3a below. The agreement loop starts with box 40, whereby for each participating network

15   device A, B, C, D as a candidate device that is represented by a candidate device

16   identifier a performs the following steps within an order.

17   A vote message comprising the candidate device identifier a, and either a first agree-value

18   Y together with the candidate value $w_a$ and the candidate justification $\pi_a$, or a second

19   agree-value N is broadcast to all participating network devices A, B, C, D, as indicated

20   with box 40. This vote messages are then received and up to $n\text{-}t$ vote messages are

21   counted, as indicated with box 50. The vote messages include either the second

22   agree-value N or the first agree-value Y, the candidate value $w_a$, and the candidate

23   justifications $\pi_a$. The candidate value $w_a$ and the candidate justifications $\pi_a$ have to satisfy

24   the predicate $Q$, otherwise they are not valid. The step in box 60 is an extension of the

25   current agreement and relates to box 30 which are described with reference to Fig. 3a and

1  3b below. As indicated in box 70, it follows a binary Byzantine agreement (BA). This

2  Byzantine agreement is validated and biased to 1 or Yes. It is used to determine whether

3  the candidate device has sent the candidate value $w_a$ and the candidate justification $\pi_a$

4  satisfying the predicate $Q$. It outputs a binary decision Y or N, which stands for Yes or

5  No, as indicated with box 80. If the decision outputs N, the protocol carries on at box 40.

6  If the decision outputs Y, a decision on the common value $v$ is performed, as indicated by

7  box 90. Thereby the common value $v$ was proposed as the candidate value $w$ and the

8  justification $p$ was proposed as the candidate justification $\pi$ of an agreed candidate device.

9  A further example embodiment relates to a constant-round protocol for multi-valued

10  Byzantine Agreement that guarantees termination within a constant expected number of

11  rounds. In the multi-valued Byzantine Agreement above, an adversary might know the

12  order in which the participating network devices A, B, C, D search for an acceptable

13  candidate device, i.e., one that has broadcast a valid proposal. Although at least one third

14  of all participating network devices A, B, C, D are guaranteed to be accepted, the

15  adversary can choose the corruptions and schedule messages such that none of them is

16  examined early in the agreement loop.

17  The remedy for this problem is to choose $\Pi$ randomly during the protocol after making

18  sure that enough participating network devices A, B, C, D already committed to their

19  votes on the candidate devices. This can be achieved in two steps. First, one round of

20  commitment exchanges is added before the agreement loop, as indicated with box 30 in

21  Fig. 2 already. Each participating network device A, B, C, D has to commit to the votes

22  that it will cast by broadcasting the identities of the $n-t$ participating network devices A,

23  B, C, D from which it has received valid echo messages. Thereby, at least consistent

24  broadcast should be used. Honest participating network devices A, B, C will later only

25  accept vote messages that are consistent with the commitments made before. The second

26  step is to determine the permutation $\Pi$ using a threshold coin-tossing scheme that outputs

27  a random, unpredictable value after enough votes are committed. Taken together, these

1    steps ensure that the fraction of participating network device A, B, C, D which are

2    guaranteed to be accepted are distributed randomly in Π, causing termination in a

3    constant expected number of rounds.

4    More precisely, the multi-valued Byzantine Agreement can be modified at the boxes 30

5    and 60 in Fig. 1, whereby the boxes 30 and 60 are described with reference to Fig. 3a and

6    3b in the following. As indicated with box 32, a commit message comprising sender

7    identities of received echo messages is broadcast. The commit messages are received, as

8    indicated with box 34. At least one cryptographic common coin is opened for selecting

9    randomly the order of the candidate device, as indicated with box 36. The opening of the

10   cryptographic common coin may comprise the use of a distributed coin-tossing protocol

11   as described above. As indicated with box 60 in Fig 3b, the vote messages are counted

12   which are consistent with the received commit messages.

13   *Atomic Broadcast*

14   Atomic broadcast guarantees a total order on messages such that honest participating

15   network devices A, B, C deliver all messages with a common tag in the same order. The

16   atomic broadcast protocol described in this invention builds directly on the multi-valued

17   Byzantine agreement.

18   The atomic broadcast protocol is efficient so that a payload message $m$ is scheduled and

19   delivered within a fixed number of steps after it is broadcast by the honest participating

20   network device A, B, C. But since the adversary may delay the sender arbitrarily and

21   deliver an a priori unbounded number of messages among the remaining honest

22   participating network devices A, B, C, the protocol can only provide such a guarantee

23   when $t+1$ honest participating network devices A, B, C become "aware" of $m$. Here the

24   definition of fairness requires that after $t+1$ honest participating network devices A, B, C

25   have broadcast some payload, it is guaranteed to be delivered within a fixed number of

26   steps. It can be interpreted as a termination condition for the broadcast of a particular

1 payload *m*. A client application might be able to satisfy this precondition through external

2 means and achieve guaranteed fair delivery in this way.

3 *Protocol for Atomic Broadcast*

4 An example protocol for atomic broadcast based on validated Byzantine agreement is

5 described in the following. Its overall structure is similar to the protocol of V. Hadzilacos

6 and S. Toueg, "Fault-tolerant broadcasts and related problems," in Distributed Systems

7 (S. J. Mullender, ed.), New York: ACM Press & Addison-Wesley, 1993 (an expanded

8 version appears as Technical Report TR94-1425, Department of Computer Science,

9 Cornell University, Ithaca NY, 1994) for the crash-fault model, but additional measures

10 to tolerate Byzantine faults are used. The atomic broadcast protocol proceeds as follows.

11 Each participating network device A, B, C, D maintains an implicit queue of not yet

12 delivered payload messages. The received messages are placed in this queue whenever

13 they are received. The protocol proceeds in asynchronous global rounds, where each

14 round comprises the following general steps:

15 - Send the current queue $q$ to all participating network devices A, B, C, D, accompanied

16 by a digital signature $\sigma$ or short signature $\sigma$.

17 - Collect the queues of $n-t$ distinct participating network devices A, B, C, D and store

18 them in a queue vector QV, store the corresponding signatures $\sigma_i$ in a signature vector

19 SV, and propose QV validated by SV for multi-valued Byzantine agreement.

20 - Perform multi-valued Byzantine agreement with validation of the queue vector QV =

21 $[q_1,..., q_n]$ by the signature vector SV = $[\sigma_1,..., \sigma_n]$ through a determined predicate Q (QV,

22 SV) which is true if and only if for at least $n-t$ distinct indices $j$, the vector element $\sigma_j$ is a

23 valid signature on a message comprising $q_j$ and $\sigma_j$ by the respective participating network

24 device.

CH920000062US1

1    - After deciding on a vector of queues DQV, deliver the union of all payload messages in

2    the decided queue vector DQV according to a deterministic order; and proceed to the next

3    round.

4    In order to ensure liveness of the atomic broadcast protocol, there are at least two ways in

5    which messages can be inserted into the queue: when a participating network device A,

6    B, C, D receives a broadcast activation message and when a participating network device

7    A, B, C, D receives the queue of another participating network device A, B, C, D

8    pertaining to the current round. If either of these two messages arrive and comprise any

9    yet undelivered payload message, and if the participating network device A, B, C, D has

10   not yet sent its own queue of the current round, then it starts the next iteration by inserting

11   the payload in its queue and sending the queue to all participating network devices A, B,

12   C, D. A more detailed description is found below with reference to Fig. 4.

13   The term $n-t$ in the atomic broadcast protocol above could be replaced by any $c$ between

14   $t+1$ and $n-t$ if the fairness condition is changed such that $n - c + 1$ participating network

15   devices A, B, C, D have to be activated instead of $t+1$.

16   Fig. 4 shows the flow of an atomic broadcast protocol 400 for an implementation. This

17   atomic broadcast protocol 400 provides a method for reliably broadcasting messages in an

18   order within the asynchronous network 1 to 5 comprising $n$ participating network devices

19   A, B, C, D. It tolerates a number $t$ of less than $n/3$ faulty participating network devices.

20   Each participating network device A, B, C, D stores a queue $q$ and a log file, hereafter

21   called log $d$, as it is known in the art. The method operates in rounds, whereby each round

22   comprises the following steps. Each participating network device A, B, C, D is

23   responsive to a message broadcast request comprising a message value $m$. If the

24   participating network device A, B, C, D receives such a broadcast request it appends the

25   message value $m$ to the queue $q$ unless the log $d$ or the queue $q$ comprises the message

26   value $m$ already, as indicated with box 410. Then, as indicated with box 420, a signature

CH920000062US1

1  $\sigma$ on a defined portion of the queue $q$ is derived. In the next step, as indicated with box

2  430, a queue message comprising the queue $q$ and the signature $\sigma$ is broadcast to all

3  participating network devices A, B, C, D. As indicated with box 440, a number $c$ of at

4  least $t+1$ queue messages comprising $c$ proposed queues $q_1$, $q_2$, $q_3$ and proposed signatures

5  $\sigma_1$, $\sigma_2$, $\sigma_3$ is then received. The proposed queues $q_1$, $q_2$, $q_3$ are stored in the queue vector

6  QV and the proposed signatures $\sigma_1$, $\sigma_2$, $\sigma_3$ are stored in the signature vector SV. This is

7  indicated in box 450. Then, as indicated with box 100, the queue vector QV is proposed

8  for a Byzantine agreement which is validated by the signature vector SV. A method for

9  achieving validated agreement on a common value is performed. Such a method can be

10  the method for achieving agreement as described with reference to Fig. 2. This method

11  takes as input a proposed value $w$, a proposed justification $\pi$, and a proposed

12  predetermined predicate $Q$. The proposed value $w$ is set to the queue vector QV, the

13  proposed justification $\pi$ is set to the signature vector SV, and the proposed predetermined

14  predicate $Q$ is set to the determined predicate Q. The determined predicate Q asserts that

15  the signature vector SV comprises $c$ valid signature entries of distinct participating

16  network devices A, B, C on entries of the queue vector QV. As indicated with box 460

17  and in response to the result of the Byzantine agreement an ordered list L of unique

18  message values out of the entries of the decided queue vector DQV is prepared. Then, the

19  unique message values are accepted in the ordered list L in the sequence of the ordered

20  list L, as indicated in box 470, and the accepted unique message values are appended to

21  the log $d$ before the next round can start, as indicated with box 480 and the arrow to the

22  start.

23  *Digital Notary Services*

24  In the following a digital notary service is described which takes advantages of the

25  methods described above. This is an example of the general state machine replication

26  technique, for which atomic broadcast protocols are useful. In general, the simplest state

27  machine is a counter. Despite this simplicity, there are a number of applications in which

1　a counter provided by a central authority is of fundamental importance. A digital notary

2　service is an example of this.

3　In its most basic form, a digital notary service receives documents from clients, assigns a

4　sequence number to each of them, and certifies this by its signature. The service provides

5　essentially a logical time stamp.

6　Assumed, the notary service is a single physical location, the documents can simply be

7　given consecutive numbers as they come in the front door and the community at large has

8　simply to trust that what happens inside that building is honest and reliable; something

9　not much different can also be done when the notary service moves to a single server on

10　the Internet. It is, however, wise to distribute such a service among $n > 3$ hosts, servers,

11　or participating network devices running different software under separate system

12　administration and in separate physical locations, in order to remove the single point of

13　failure and to provide a more robust service. A reasonable model to apply for the

14　distributed servers is one of participating network devices A, D, C, D communicating

15　over the asynchronous network 1 to 5, which is presumed to be in the hands of the

16　adversary D, who is also presumed to have corrupted up to $t < n/3$ of the participating

17　network devices or servers. Several security concerns have to be satisfied if the

18　distributed notary service is to continue to function:

19　• The documents should be processed atomically, i.e., despite asynchronous

20　communication, each application receives a unique sequence number, and all honest

21　participating network devices A, B, C have to agree on the number assigned to a

22　given document.

23　• The resulting sequence numbers should be signed in such a way that the client can be

24　confident that her number is provably legitimate, and not merely the result of

25　deceptive messages from the adversary.

CH920000062US1

1      These requirements can be met by using an atomic broadcast protocol for broadcasting

2      the documents in an order. The notary application has to take some additional steps to

3      produce the desired signatures. In particular:

4      - During initialization, a dealer creates keys for an $(n, t+1)$-threshold signature scheme $S$

5      and gives them to the participating network devices A, B, C, D, hereafter also referred to

6      as notary servers. Within the initialization phase, each notary server initializes a counter

7      *seqnum* to zero.

8      - The notary servers wait for the delivery of a message by the atomic broadcast protocol,

9      whose payload $m$ comprises the document to be notarized and the identity of the client.

10      - When such a message is delivered, each notary server increments the counter *seqnum*,

11      creates an $S$ -signature share $\sigma$ on the message $(m, seqnum)$, and transmits $(seqnum, \sigma)$ to

12      the client indicated in the delivered message.

13      The client uses this service by

14      - sending the document $m$ in a message to all notary servers, requesting atomic broadcast

15      delivery of this message among the notary servers,

16      - waiting for at least $t+1$ distinct notary servers to reply with a sequence number and valid

17      $S$-signature share and then

18      - assembling the shares into a $S$-signature on the message $(m, seqnum)$.

19      One inelegant aspect of the notary service is that it seems to put some of the onus of

20      dealing with the distributed service on the individual clients, in that the client has send

21      initially a message to all notary servers and to assemble herself the signature shares at the

22      end of the process. This could be avoided, at the cost of an additional round of

23      communication among the notary servers, by having each one broadcast its $S$-signature

CH920000062US1

1     share to its peers and then wait for enough such shares to come back from other notary

2     servers so that it could assemble the signature itself. However, in practice, the sending of

3     the client's outgoing message and the assembly of a signature from the incoming server

4     messages would be handled by routines in the communications software library, so these

5     internals would not be the concern of the client application at all.

6     *Other Distributed Services*

7     Several other distributed services with a similar need of a synchronized and signed

8     sequence number can all fall under the rubric of a "digital bidding service". One can

9     imagine a valuable item for which the order of arrival of the bids as well as the details of

10     the bids themselves are used to allocate the good or goods to (some of) the clients.

11     Stocks, where both the order in which the offers are made as well as the offered purchase

12     price, or government contracts, where again the priority of the bid as well its specific

13     terms are used to assign the contract, are both items which would benefit from such a

14     service. In individual applications, the details of the small modifications to the atomic

15     broadcast protocol which would be necessary may vary, such as particular requirements

16     of client authentication or creation of a mechanism to terminate a bidding process, but the

17     general outline would be very similar to the above example.

18     *Hybrid adversary structures*

19     Instead of a fixed threshold of $t$ out of $n$ corruptions, it is possible to gain more flexibility

20     by reflecting real world structures.

21     For example, an adversary could be able to control all participating network devices with

22     a certain operating system, or he might bribe one system administrator to get access to all

23     participating network devices at a specific site. Adversary structures cope with such an

24     attack scheme.

CH920000062US1

1    In general, to define an adversary structure $T$, one has to define every coalition of parties

2    whose corruption the system should tolerate, e.g., a coalition of all participating network

3    devices with the same operating system. The set of all those sets then is the adversary

4    structure $T$.

5    In the method for achieving agreement among $n$ participating network devices and the

6    method for reliably broadcasting messages in an order, several types of failures can occur

7    simultaneously. For example, it could differ between crash failures CF, Byzantine failures

8    BF, and link failures LF. This allows for a higher number overall number of failures to be

9    tolerated.

10    The present invention can be realized in hardware, software, or a combination of

11    hardware and software. A visualization tool according to the present invention can be

12    realized in a centralized fashion in one computer system, or in a distributed fashion where

13    different elements are spread across several interconnected computer systems. Any kind

14    of computer system - or other apparatus adapted for carrying out the methods and/or

15    functions described herein - is suitable. A typical combination of hardware and software

16    could be a general purpose computer system with a computer program that, when being

17    loaded and executed, controls the computer system such that it carries out the methods

18    described herein. The present invention can also be embedded in a computer program

19    product, which comprises all the features enabling the implementation of the methods

20    described herein, and which - when loaded in a computer system - is able to carry out

21    these methods.

22    Computer program means or computer program in the present context include any

23    expression, in any language, code or notation, of a set of instructions intended to cause a

24    system having an information processing capability to perform a  particular function

CH920000062US1

1    either directly or after conversion to another language, code or notation, and/or

2    reproduction in a different material form.

3    Thus the invention includes an article of manufacture which comprises a computer usable

4    medium having computer readable program code means embodied therein for causing a

5    function described above. The computer readable program code means in the article of

6    manufacture comprises computer readable program code means for causing a computer to

7    effect the steps of a method of this invention. Similarly, the present invention may be

8    implemented as a computer program product comprising a computer usable medium

9    having computer readable program code means embodied therein for causing a a function

10    described above. The computer readable program code means in the computer program

11    product comprising computer readable program code means for causing a computer to

12    effect one or more functions of this invention. Furthermore, the present invention may be

13    implemented as a program storage device readable by machine, tangibly embodying a

14    program of instructions executable by the machine to perform method steps for causing

15    one or more functions of this invention.

16    It is noted that the foregoing has outlined some of the more pertinent objects and

17    embodiments of the present invention. This invention may be used for many

18    applications. Thus, although the description is made for particular arrangements and

19    methods, the intent and concept of the invention is suitable and applicable to other

20    arrangements and applications. It will be clear to those skilled in the art that

21    modifications to the disclosed embodiments can be effected without departing from the

22    spirit and scope of the invention. The described embodiments ought to be construed to

23    be merely illustrative of some of the more prominent features and applications of the

24    invention. Other beneficial results can be realized by applying the disclosed invention in

25    a different manner or modifying the invention in ways known to those familiar with the

26    art.